

EV369763178

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

**DATA OVERLAY, SELF-ORGANIZED METADATA OVERLAY,
AND ASSOCIATED METHODS**

Inventors:

Zheng Zhang

and

Shu-Ming Shi

ATTORNEY'S DOCKET NO. MS1-1959US

TECHNICAL FIELD

This invention relates to a distributed data structure and to a technique for utilizing the data structure to interact with a peer-to-peer system.

BACKGROUND

Peer-to-peer (P2P) systems employ a network which connects participating machines having equal or similar capabilities and responsibilities. These systems perform tasks without the coordination of a conventional server (or with minimal set-up coordination by a server). For instance, Fig. 1 shows a high-level depiction of a P2P system 100. The system 100 includes a collection of peer entities (102-112) having equal or similar capabilities and responsibilities. In one example, the peer entities (102-112) may correspond to independent personal computer devices coupled together via an Internet or intranet. The peer entities (102-112) can directly transfer files or other information between themselves (as indicated by exemplary communication path 114) without the aid of a server. A general introduction to P2P systems can be found in D. S. Milojevic, V. Kalogeraki, R. Lukose, K. Nagaraja, J. Pruyne, B. Richard, S. Rollins, and Z. Xu., "Peer-To-Peer Computing," Technical Report HPL-2002-57, HP Lab, 2002.

P2P systems commonly use a distributed hash table (DHT) to facilitate the storage and retrieval of objects from peer entities participating in the systems. As the name suggests, a distributed hash table (DHT) refers to a hash table that is distributed over plural locations, such as distributed over plural stores associated with different computer devices. A distributed hash table specifies a plurality of DHT nodes having respective assigned IDs. The DHT nodes collectively define an abstract DHT logical space. An object can be inserted into or retrieved from this DHT logical space by subjecting this object to a hashing function to produce a key. This key is then used to locate a particular

1 target node ID in the DHT logical space that will receive the object or from which the
2 object can be retrieved. That is, each DHT node is associated with a range of keys;
3 objects are added to or retrieved from a particular DHT node depending on whether their
4 key falls within the range of keys associated with that particular DHT node. Unlike non-
5 distributed hash table implementations, DHT nodes can freely join and leave the DHT
6 logical space (e.g., corresponding to computer devices joining and leaving the P2P
7 system, respectively), so functionality must be provided to address these events.

8 A variety of DHT strategies have been developed to manage the storage and
9 retrieval of objects in a P2P system. Fig. 2 shows a Content Addressable Network (CAN)
10 strategy, e.g., as described in S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S.
11 Shenker, "A Scalable Content-Addressable Network," ACM SigComm 2001, San Diego,
12 CA, USA, Aug. 2001. This strategy models the DHT logical space as a D-dimensional
13 Cartesian space 200. The CAN strategy partitions the space 200 as nodes join the DHT
14 space 200. For instance, when node n1 joins, the CAN strategy allocates the entire space
15 200 to this node. When node n2 joins, the CAN strategy divides the space 200 into two
16 halves and allocates each half to nodes n1 and n2, respectively. When node n3 joins, the
17 CAN strategy divides the left half into upper and lower quarters, assigning the upper
18 quarter to node n2 and the lower quarter to node n3. And when node n4 joins, the CAN
19 strategy divides the lower quarter into a left eighth (which is assigned to node n3) and a
20 right eighth (which is assigned to node n4). This procedure is repeated as many times as
21 necessary to dynamically account for nodes being adding and removed. The resultant
22 partitions define logical spaces used to insert and retrieve objects into and from the
23 distributed hash table. A node can be said to "own" the objects that map to its space.

24 Fig. 3 shows another strategy referred to as CHORD (e.g., as described in I.
25 Stoica, R. Morris, D. Karger, M.F. Kaashoek, and H. Balakrishnan, "Chord: a Scalable

1 Peer-To-Peer Lookup Service for Internet Applications,” ACM SigComm 2001, San
2 Diego, CA, USA, Aug. 2001. In this strategy, the DHT logical space is structured as
3 circular space 300. DHT nodes are assigned IDs and added to the circular DHT logical
4 space 300 based of their assigned IDs. For instance, exemplary DHT nodes n1, n2, n3,
5 n4, and n5 shown in Fig. 3 have assigned IDs that govern their “placement” on the
6 circular DHT logical space 300. As in the case of Fig. 2, the DHT nodes partition the
7 DHT logical space 300 as they are added, defining multiple subspaces or zones. These
8 zones define the objects that each node “owns.” For instance, to insert an object into a
9 distributed hash table that is governed by the DHT strategy shown in Fig. 3, the object is
10 subjected to a hashing function to produce a key. The object is then stored at the DHT
11 node have a zone assigned to that key (e.g., at the DHT node which encompasses a range
12 of keys that include the object’s key). In both the cases of Fig. 2 and Fig. 3, a variety of
13 lookup strategies can be used to quickly find a particular node in the P2P system. In
14 general, the lookup strategies involve making several “hops” in the DHT logical space to
15 narrow in on the desired target DHT node. Various mechanisms are commonly provided
16 to expedite this search. For instance, each DHT node in the CHORD strategy stores the
17 IDs of a set of other DHT nodes. These other IDs can increase in exponential fashion,
18 establishing so-called “fingers” that probe out into the logical space 300. This allows the
19 lookup procedure to quickly locate a desired DHT node with a small number of hops.

20 Figs. 2 and 3 provide merely a high level overview of two exemplary known DHT
21 routing strategies. There are many other strategies. For instance, another popular routing
22 strategy is the PASTRY routing strategy, as described in A. Rowstron and P. Druschel,
23 “Pastry: Scalable, Distributed Object Location and Routing for Large-Scale Peer-To-Peer
24 Systems,” 18th FIFP/ACM International Conference on Distributed Systems Platforms
25 (Middleware), Heidelberg, Germany, Nov. 2001.

1 P2P systems offer many benefits over conventional client-server strategies. For
2 instance, P2P systems have the ability to automatically and freely expand and contract
3 without central coordination. But this lack of supervisory coordination also poses various
4 challenges. For instance, it may be desirable to have the P2P system act in concert to
5 perform some global function. For instance, in various instances, it may be desirable to
6 collect data from the participants of the P2P system. Or it may be desirable to
7 disseminate information to the participants in the P2P system. With a client-server
8 approach, a server can simply poll its clients to collect information from its clients, or
9 broadcast information to its clients to disseminate information to its clients. But data
10 gathering and dissemination becomes more problematic in a P2P system because it is
11 formed by a loose alliance of interconnected peers that can freely come and go. Adding
12 centralized conventional reporting functionality may have the effect of complicating the
13 P2P system, and thus reducing its flexibility and utility.

14 There is accordingly an exemplary need in the art for an efficient strategy for
15 interacting with a P2P DHT that will allow, for instance, for the gathering of data from its
16 participants and the dissemination of information to its participants.

17 18 **SUMMARY**

19 According to one exemplary implementation, a method is described for building a
20 data overlay. The method includes providing a distributed hash table (DHT) that governs
21 the insertion and retrieval of objects into and from a peer-to-peer system, wherein the
22 distributed hash table includes a logical space including a plurality of DHT nodes having
23 an associated plurality of DHT zones. The method also includes building the data
24 overlay as a data structure on top of the logical space of the distributed hash table by
25

1 associating objects in the data structure with the DHT nodes, and by establishing links
2 between the objects in the data structure.

3 According to another exemplary implementation, the above-described data
4 overlay can have a topology of a tree, the tree having a plurality of tree nodes associated
5 with respective DHT nodes, wherein each tree node has a respective tree node zone
6 associated therewith which corresponds to a part of the logical space of the distributed
7 hash table.

8 According to another exemplary implementation, a method is described for
9 passing data through a data overlay. As in the previous case, the method includes
10 providing a distributed hash table (DHT) that governs the insertion and retrieval of
11 objects into and from a peer-to-peer system, wherein the distributed hash table includes a
12 logical space including a plurality of DHT nodes having a plurality of associated DHT
13 zones. The method further includes building a data overlay as a data structure on top of
14 the logical space of the distributed hash table by associating objects in the data structure
15 with the DHT nodes, and by establishing links between the objects in the data structure,
16 wherein the data overlay defines a plurality of interconnected nodes. In this
17 implementation, the method further includes routing data through the data overlay by
18 passing the data through its interconnected nodes.

19 According to another exemplary implementation, the above-described data
20 overlay used for routing has a topology of a tree, the tree having a plurality of tree nodes
21 associated with respective DHT nodes. The routing of data through the data overlay can
22 include gathering data from DHT nodes and passing the data up through the tree nodes to
23 a root node of the tree. The routing of data through the data overlay can also include
24 disseminating data from the root node of the tree, through the tree nodes, to the DHT
25 nodes.

Additional implementations and features will be described in the following.

BRIEF DESCRIPTION OF THE DRAWINGS

Fig. 1 shows a conventional peer-to-peer (P2P) system.

Fig. 2 shows a conventional CAN routing strategy.

Fig. 3 shows a conventional CHORD routing strategy.

Fig. 4 shows a conventional technique for linking two objects of a data structure in the context of a local machine environment.

Fig. 5 shows an exemplary technique for linking two objects of a data structure in a P2P distributed hash table (DHT) environment, where the two objects are associated with two different nodes in the P2P DHT environment. The linking technique forms the basis of a data overlay placed “on top” of the DHT.

Fig. 6 shows an exemplary tree structure constructed using the concept of a data overlay depicted in Fig. 5. The tree structure is referred to a Self-Organized Metadata Overlay (SOMO).

Fig. 7 shows an exemplary procedure for building the SOMO tree structure of Fig. 6.

Fig. 8 shows an exemplary application of the SOMO tree structure of Fig. 6 to the collection information from the participants of the P2P system.

Fig. 9 shows an exemplary application of the SOMO tree structure of Fig. 6 to the dissemination of information to the participants of the P2P system.

Fig. 10 shows an exemplary computer used to implement a participant of a P2P system, where the P2P system includes a data overlay built on top of its DHT.

The same numbers are used throughout the disclosure and figures to reference like components and features. Series 100 numbers refer to features originally found in Fig. 1,

1 series 200 numbers refer to features originally found in Fig. 2, series 300 numbers refer
2 to features originally found in Fig. 3, and so on.

3 4 **DETAILED DESCRIPTION**

5 The strategies described herein pertain to a data structure built “on top” of a
6 distributed hash table (DHT) used in a peer-to-peer (P2P) system. The term peer-to-peer
7 (P2P) system can describe any interconnection of participants in which the participants
8 can directly interact with others, such as the interconnection network 100 shown in Fig.
9 1. In one implementation, the P2P system does not require the assistance of any server-
10 type entities. The participants can include any kind of entity, including personal
11 computers, laptop computers, personal digital assistants, application-specific computing
12 devices, and so on. The participants can communicate with each other via any
13 combination of routing infrastructure, such as hardwired and/or wireless communication
14 routing mechanisms, various routers, gateways, etc. Further, the participants can
15 communicate with each other through any combination of network protocols, such as
16 TCP/IP (e.g., as provided by the Internet or an intranet).

17 More generally, any of the functions described herein can be implemented using
18 software, firmware (e.g., fixed logic circuitry), manual processing, or a combination of
19 these implementations. The term “logic” or “module” as used herein generally represents
20 software, firmware, or a combination of software and firmware. For instance, in the case
21 of a software implementation, the term “logic” or “module” represents program code that
22 performs specified tasks when executed on a processing device or devices (e.g., CPU or
23 CPUs). The program code can be stored in one or more computer readable memory
24 devices.

1 This disclosure includes the following: Section A describes a general data overlay
2 structure that can be built “on top” of a P2P DHT; Section B describes an exemplary tree
3 structure built using the concept of a data overlay set forth in Section A; Section C
4 describes an exemplary method for building the tree structure and for subsequently
5 applying it to the gathering and dissemination of information in a P2P system; and
6 Section D describes an exemplary P2P participant that can be used in the type of P2P
7 DHT system described in Sections A-C.

8 9 **A. The Data Overlay**

10 The novel concept of a “data overlay” is described here. A data overlay is a data
11 structure comprised of objects. The data structure is implemented “on top” of a
12 distributed hash table. By way of background, a distributed hash table (DHT) provides a
13 technique for inserting objects into and retrieving objects from a distributed store
14 provided by a P2P system. It performs this task by defining a collection of DHT nodes
15 within a logical DHT space. That is, the DHT technique assigns each DHT node to a
16 predetermined portion of the DHT logical space, referred to as the “zone” of the DHT
17 node. For example, in the CHORD technique, a particular DHT node’s zone can be
18 interpreted as the span defined between that particular DHT node and its adjacent node in
19 a circular DHT logical space (e.g., as shown in Fig. 3). An object is stored by hashing it
20 to produce a key, and then using this key to associate the object with a particular node ID
21 in the DHT logical space. The object is retrieved from the DHT logical space in a related
22 manner. The associated zones ultimately map into real machines (e.g., computer devices
23 and associated file storage systems), although there need not be a one to one relationship
24 between nodes and machines.
25

1 The data overlay is implemented “on top” of the DHT in the sense that its objects
2 are associated with nodes in the DHT logical space. Further, an application traverses (or
3 routes) from one object to another in the data overlay’s data structure using the
4 underlying protocols and services of the P2P DHT. More specifically, for frame of
5 reference, consider the conventional case of Fig. 4 of a single machine environment 402.
6 In this environment 402, a data structure includes two objects, *a* 404 and *b* 406,
7 implemented in the storage provided by a single machine. An object broadly represents
8 any unit of any type of information; in a conventional case, for instance, an object can
9 correspond to a database record, e.g., a document. In the example of Fig. 4, object *a* 404
10 contains a pointer 408 that references object *b* 406.

11 In contrast, Fig. 5 shows the implementation of a data overlay in the context of a
12 P2P DHT environment 502. In this environment 502, since the objects are built “on top”
13 of the DHT node framework already provided by the DHT, individual nodes in the DHT
14 logical space “host” the objects in the data overlay. For example, DHT node *x* 504 hosts
15 object *a* 506, and DHT node *y* 508 hosts object *b* 510. In this example, object *a* 506
16 references object *b* 510. Generally, object *a* 506 can link to object *b* 510 by storing the
17 key that is used to access object *b* 510. This key is established when object *b* 510 is
18 created. In the case of Fig. 5, however, the referencing scheme includes two fields. A
19 first field 512 contains a hardwired address that points from object *a* 506 to object *b* 510.
20 This field is referred to as *a.foo.key*. A second field 514 contains a soft-state reference
21 that identifies the last known DHT node (e.g., node *y* 508) that hosts object *b* 510. This
22 field is referred to as *a.foo.host*. The second field 514 thus serves as a routing shortcut to
23 access object *b* 510.

24 Because the nodes of the data overlay can be dispersed over plural DHT nodes,
25 the data overlay itself can be regarded as a distributed data structure. Although the data

1 structure is distributed, it may be desirable to store it in such a manner that its objects are
2 not unduly widely geographically dispersed. This can be achieved by generating the keys
3 of a 506 and b 510 so that they are close to each other. This makes it more likely that the
4 P2P DHT system will associate these keys with the same node in the P2P system or in
5 closely related nodes in the P2P DHT system.

6 The data overlay also provides a collection of primitives used to manipulate
7 pointers and objects in its data structure. More specifically, these primitives include a
8 procedure (setref) for establishing a reference from object a to another object b , a
9 procedure (deref) for returning an object pointed to by object a , and a procedure for
10 deleting (delete) an object pointed to by object a .

11 Because the data overlay is implemented on top of the DHT system, its primitives
12 use the DHT's services. For example, the primitives can use a DHT_insert service for
13 inserting an object into the DHT logical space. The primitives can use a DHT_lookup
14 service for using a predetermined DHT routing procedure to find an object based on its
15 key in the DHT logical space (such as the exponential finger lookup structure used by
16 CHORD). And the primitives can also use a DHT_direct procedure for directly accessing
17 an object if the DHT node that stores the object is known in advance. In other words,
18 DHT_direct bypasses the normal DHT_lookup routing procedure and directly seeks the
19 node that hosts the object given its key. Both DHT_lookup and DHT_insert will, as a
20 side effect, return the DHT node in the DHT that currently hosts the target object.

21 More specifically, an exemplary procedure for establishing a pointer from one
22 object to another is provided in the following exemplary pseudo-code:
23
24
25

Pseudo-Code Excerpt No. 1: Exemplary setref Primitive

```
1      setref(a.foo, b) {          // initially a.foo==null; b is the object
2
3          // to which a.foo will point to
4
5          a.foo.key=b.key
6
7          a.foo.host=DHT_insert(b.key, b)
8
9      }
```

10 In this procedure, the first field in the referencing scheme, a.foo.key, is assigned the key
11 value of the object which it points to (that is, b.key). The second field in the referencing
12 scheme, a.foo.host, is assigned the value returned by the DHT_insert service.

13 As exemplary procedure for returning an object pointed to by a.foo is provided in
14 the following exemplary pseudo-code:

Pseudo-Code Excerpt No. 2: Exemplary deref Primitive

```
14      deref(a.foo) {            // return the object pointed to by a.foo
15
16          if (a.foo!=null) {
17
18              obj=DHT_direct(a.foo.host, a.foo.key)
19
20              if obj==null { // object has moved
21
22                  obj=DHT_lookup(a.foo.key)
23
24                  a.foo.host = node returned
25
26              }
27
28          return obj
29
30          else return "non-existent"
31
32      }
```

1 This procedure attempts to directly retrieve the object *b* using DHT_direct based on the
2 a.foo.host field and the a.foo.key field. However, if this is not successful, the procedure
3 attempts to retrieve the object *b* using the DHT_lookup service based on the a.foo.key.
4 The DHT_lookup service also will return the node that hosts the object *b*, which is
5 assigned to the field a.foo.host.

6 Finally, an exemplary procedure for deleting an object pointed to by a.foo is
7 provided in the following exemplary pseudo-code:

8
9 Pseudo-Code Excerpt No. 3: Exemplary delete Primitive

```
10 delete(a.foo) {           // delete the object pointed to by a.foo
11     DHT_delete(a.foo.key)
12     a.foo=null
13 }
```

14
15 This procedure uses the DHT_delete service of the underlying DHT to delete the object *b*
16 based on the field a.foo.key.

17 The data overlay can be implemented using its underlying DHT service by
18 modifying whatever library routines are used to create objects so that these routines also
19 establish the pointers described above as attributes of the objects. The library routines
20 should also be modified to accommodate the primitives described above for setting a
21 reference, returning an object pointed to by a reference, and deleting an object pointed to
22 by a reference.

23 There are a number of advantages to building the data overlay on top of a DHT.
24 For example, the DHT is designed to self-organize as DHT nodes are added to and
25 deleted from the DHT logical space (related to real machines joining and leaving the P2P

1 system, respectively). The DHT is also designed to automatically “heal” itself in
2 response to DHT nodes being added to and deleted from the DHT logical space (such as
3 by reestablishing links between nodes, transferring objects between nodes, etc.). By
4 virtue of being implemented on top of the DHT, the data overlay can also adopt the
5 features of self-organizing and self-healing. More specifically, the data overlay can be
6 configured such that it is self-organizing and self-healing on the same scale as the
7 underlying DHT.

8 Further, various applications can be ported to run on top of a P2P DHT, giving
9 these applications the illusion of an infinite storage space (e.g., giving the impression of a
10 single resource pool having a large size encompassing the nodes of the DHT logical
11 space). This storage space can broadly include memory heaps of machines that are
12 participating in the P2P DHT system. The host routing shortcut (e.g., a.foo.host) makes
13 the performance of applications utilizing the data overlay independent of the underlying
14 DHT system.

15 16 **B. The SOMO Tree Structure; an Instance of the Data Overlay**

17 The above data overlay provides a framework for building an arbitrary data
18 structure on top of a DHT. The data structure includes a plurality of objects which
19 constitute nodes in the data structure. This data structure can assume any kind of
20 topology by linking the nodes together in different ways. Further, the data structure can
21 implement different functions depending on the operations assigned to its individual
22 nodes. The following section described an exemplary instance of the data overlay data
23 overlay referred to as a Self-Organized Metadata Overlay, or “SOMO” for brevity.

24 The SOMO data structure is constructed to assume the topology of a tree
25 structure. The SOMO tree structure has a root node. The root node can have one or more

1 children, which, in turn, can have their own respective children. The terminal nodes of
2 the SOMO tree structure are referred to as leaf nodes. The leaf nodes are associated with
3 respective DHT nodes in the DHT logical space of the P2P DHT system.

4 As will be described in greater detail in Section C below, one function of the
5 SOMO tree structure is to extract metadata from the DHT nodes (which ultimately
6 involves extracting data from the machines that implement the P2P system) and to pass
7 this metadata up through the SOMO tree to the root node of the SOMO tree structure. An
8 application can then read this metadata and perform some action on the basis of this
9 metadata. (Metadata generally refers to any kind of information associated with the
10 operations being performed in the P2P system, such as information regarding the
11 performance of machines that comprise the P2P system). The SOMO tree structure can
12 also be used to disseminate information from the root node of the SOMO tree structure
13 down to the DHT nodes and associated machines within the P2P system. Thus, generally
14 speaking, the SOMO tree structure can serve the role of data gathering (e.g., aggregation)
15 and data broadcast.

16 Fig. 6 illustrates an exemplary SOMO tree structure 602 that is built on top of an
17 underlying DHT logical space 604. The DHT logical space 604 is partitioned into a
18 number of zones, such as exemplary zone 606 and exemplary zone 608. Each zone
19 includes a DHT node associated therewith, such as exemplary DHT node 610. The DHT
20 can partition the DHT logical space 604 into zones according to any technique, such as
21 exemplary techniques provides by the CAN partitioning scheme, CHORD partitioning
22 scheme, PASTRY partitioning scheme, or any other kind of DHT partitioning scheme.
23 For example, using the CHORD partitioning scheme, the DHT logical space 604 can be
24 defined as a ring having nodes dispersed at various locations around it, and the zones can
25 correspond to the spans that separate neighboring adjacent DHT nodes on the ring.

1 The SOMO tree structure 602 includes one or more nodes that are referred to here
2 as “SOMO nodes” to distinguish them from DHT nodes. Each SOMO node is
3 represented by symbol s . The exemplary SOMO tree structure 602 shown in Fig. 6
4 includes SOMO nodes s 612-626. The nodes s 612-626 form an inverted tree shape.
5 Namely, a root node 612 branches off into child node 614 and child node 616. These
6 child nodes can have their own respective child nodes; for example, child node 614
7 includes child node 618 and child node 620. Although the full structure of the exemplary
8 SOMO tree structure 602 is abbreviated in Fig. 6 to facilitate illustration and discussion,
9 the SOMO tree structure 602 ultimately terminates in leaf nodes (e.g., leaf nodes 622,
10 624, 626) planted in corresponding DHT nodes in the DHT logical space 604. In general,
11 the links between the SOMO nodes in the SOMO tree structure 602 are illustrated in Fig.
12 6 by dotted lines that connect the SOMO nodes together; these links can be implemented
13 using the referencing scheme described in Section A above.

14 Each SOMO node s has a zone associated therewith. For example, the root
15 SOMO node 612 includes a zone 628 that spans the entire DHT logical space 604. Child
16 node 616 includes a zone 630 which spans one half of the root node 612’s zone 628.
17 Another child node 620 that is deeper in the SOMO tree structure 602 has a zone 632 that
18 is one quarter of the root node 612’s zone 628. Accordingly, successive nodes s added to
19 the hierarchy of the SOMO tree structure 602 result in progressively finer partitioning of
20 the root node 612’s zone 628. Also, the hierarchy of the SOMO tree structure 602 grows
21 “taller” for those regions of the DHT logical space 604 that exhibit finer (that is, denser)
22 partitioning of the space 604. In general, Fig. 6 represents the zones associated with
23 individual SOMO nodes by horizontal arrows that span the length the SOMO nodes’
24 respective zones. A DHT node that hosts a particular SOMO node s is expressed as
25 $\text{DHT_host}(s)$.

1 An exemplary data structure that can be used to implement each SOMO node s is
2 given by the following code excerpt:

3
4 Pseudo-Code Excerpt No. 4: Exemplary SOMO Node Data Structure

```
5       struct SOMO_node {  
6             string key  
7             struct SOMO_node *child[1..k]  
8             DHT_zone_type Z  
9             SOMO_op op  
10            Report_type report  
11       }
```

12
13 The member *report* in the above data structure indicates the type of report that a
14 particular SOMO node is entrusted to handle, such as a sorted list. The member Z in the
15 above data structure identifies the zone which this particular SOMO node's *report*
16 member covers. That is, for example, the zone Z for SOMO node 620 corresponds to
17 zone 632 shown in Fig. 6. The *key* member in the above data structure refers to a key
18 assigned to this particular SOMO node. This *key* can be generated by applying a
19 deterministic function to the SOMO node's zone Z . For instance, the *key* can correspond
20 to the center of the zone Z , or the result of a hashing function applied to the zone Z , such
21 as a hash of Z 's defining coordinates. For example, assume that a zone is defined by the
22 1D logical span of 0 to 1, e.g., $[0, 1]$; in this case, the *key* can correspond to the center
23 coordinates of this zone [i.e., 0.5], or some hash of these coordinates. By virtue of this
24 key referencing scheme, a SOMO node s will be hosted by a DHT node that covers $s.key$
25 (e.g., the center of $s.Z$). This allows a SOMO node to be retrieved deterministically as

1 long as an application knows its region. This provision is particularly useful when
2 querying system status in a given key-space range.

3 The above data structure also identifies the children of the particular SOMO node,
4 e.g., by storing pointers to these children. As describe above, the addition of children has
5 the effect of diving the SOMO node's zone Z by a factor of k . Each child is allocated to
6 one of these zone partitions. That is, a SOMO node's i -th child will cover the i -th
7 fraction of region $s.Z$. This recursion continues until a termination condition is met (to be
8 discussed below). Since a DHT node will "own" a piece of the DHT logical space 604, a
9 leaf SOMO node s is therefore guaranteed to be "planted" in it.

10 Finally, the member *op* defines an operation that the particular SOMO node can
11 perform on information that passes through it (in either a data gathering or data
12 dissemination mode.). For example, the *op* can specify that a merge-sort operation is to
13 be performed in the course of collecting information using the SOMO tree structure 602.
14 By virtue of the inclusion of the *op* member, the SOMO tree structure 602 can execute
15 any functionality in a distributed and parallel manner. Thus, the SOMO tree structure
16 602 can also be viewed as a mechanism for providing a distributed parallel processing
17 framework to implement any kind of functionality.

18 19 **C. Method for Building and Applying the SOMO Tree Structure**

20 *C.1. Building the SOMO Tree Structure*

21 As mentioned above, the data overlay can grow and shrink as a function of
22 dynamic and unsupervised modifications made in the underlying DHT. Since the SOMO
23 tree structure 602 is an instance of the data overlay, this means that the SOMO tree
24 structure 602 also has the ability to grow and shrink in response to modifications made to
25 the underlying DHT. Also, the SOMO tree structure, like its underlying DHT, has the

1 ability to heal itself to counteract modifications in the underlying DHT. This subsection
2 describes the manner in which the SOMO tree structure 602 evolves in response to
3 changes in its underlying DHT.

4 In one implementation, the SOMO tree structure 602 is grown by first adding the
5 root node 628 corresponding to an initial DHT node added to the DHT logical space 604,
6 and then successively adding additional SOMO nodes s as additional DHT nodes are
7 added to the DHT logical space 604. In this manner, the SOMO tree structure 602 grows
8 “taller” as the DHT logical space 604 becomes more finely partitioned.

9 Fig. 7 shows a high level depiction of a procedure 700 used to determine whether
10 child nodes should be added to a particular SOMO node s . In step 702, it is determined
11 whether the SOMO node’s responsible zone $s.Z$ is smaller than or equal to that of its
12 hosting DHT node’s zone. If step 702 is answered in the affirmative, then this SOMO
13 node already corresponds to a leaf planted in the correct DHT node, and there is no need
14 to add a child to this particular SOMO node. This outcome is represented by step 704 of
15 Fig. 7. Alternatively, if step 702 is answered in the negative, this means that the SOMO
16 node’s zone $s.Z$ encompasses partitions in the DHT logical space 604 that are not
17 currently accounted for by the SOMO tree structure 602. Since the SOMO tree structure
18 602 needs to interact with these DHT partitions, the procedure 700 responds by adding
19 another SOMO node s to the SOMO tree structure 602. The SOMO node s has the data
20 structure defined above.

21 The procedure shown in Fig. 7 refers to decisions and actions made with respect
22 to an individual SOMO node s in the SOMO tree structure 602. These operations are
23 repeated for each individual SOMO node in the SOMO tree structure 602. Furthermore,
24 these operations can be repeated at periodic intervals. By virtue of this procedure, the
25 SOMO tree structure 602 automatically adapts to changes in the DHT logic space 604 by

adding SOMO child nodes. As described above, the SOMO tree structure 602 will be “taller” in logical space regions where DHT nodes are denser (that is, more finely partitioned). Nodes can also be removed from the SOMO tree structure 602 to react to the removal of DHT nodes from the DHT logical space 604.

The following pseudo-code provides one exemplary routine (e.g., SOMO_grow) to add SOMO nodes to the SOMO tree structure 602:

Pseudo-Code Excerpt No. 5: Exemplary SOMO_Grow Procedure

```

SOMO_grow(SOMO_node s) {
    // check if it is necessary to add any children
    if (s.Z  $\subseteq$  DHT_host(s).Z) return
    for i=1 to k
        if (s.child[i]==NULL && the i-th sub-space of s.Z  $\not\subseteq$  host(s).Z) {
            t = new(type SOMO_node)
            t.Z = the i-th sub-space of s.Z
            t.key = SOMO_loc(t.Z)
            setref(s.child[i], t) // inject into DHT
        }
    }
}

```

This procedure involves initializing a SOMO node object and its appropriate fields, and then calling the setref primitive (defined above) to install the pointer (that links this SOMO node to the data structure); this last step involves invoking a DHT operation. As described above, the SOMO_grow procedure also involves deterministically calculating a SOMO node’s *key* given the region it covers.

1 The SOMO_loc procedure used in the above-defined routine can be implemented
2 using the following pseudo-code:

4 Pseudo-Code Excerpt No. 6: exemplary SOMO_loc procedure

```
5       SOMO_loc(DHT_zone_type Z) {  
6             return center of Z  
7             // optionally  
8             // return hash_of (Z)  
9       }
```

11 As described above, the above-described SOMO_grow procedure is performed in
12 a top down fashion, and is executed periodically. In another implementation, a bottom-up
13 protocol can be employed. When nodes are removed from the DHT logical space 604,
14 the SOMO tree structure 602 will prune itself accordingly by deleting redundant child
15 SOMO nodes.

16 For an N -node system where nodes populate the total logical space evenly, there
17 will be $2N$ SOMO nodes when the SOMO fan-out k is 2. The crash of a DHT node will
18 take away the SOMO nodes that this DHT node is hosting. However, the crashing DHT
19 node's zone will be taken over by another DHT node after repair. Consequently, the
20 periodic checking of all child SOMO nodes ensures that the SOMO tree structure 602 can
21 be completely reconstructed in $O(\log_k N)$ time (where the conventional "Big O" notation
22 is used here to express algorithm performance). Because the SOMO root node is always
23 hosted by the DHT node that owns one deterministic point of the total DHT logical space
24 604, that DHT node ensures the existence of the SOMO root node and invokes the
25 SOMO_grow routine on the SOMO root node.

1 In one case, the result of the SOMO_grow procedure is the exemplary SOMO tree
2 structure 602 shown in Fig. 6. SOMO nodes are scattered among DHT nodes. The
3 SOMO tree structure 602 can be physically implemented in distributed fashion in the
4 stores of individual machines or other infrastructure that comprise the P2P system.

6 *C.2. Applying the SOMO Tree Structure*

7 As described above, one exemplary use of the SOMO tree structure 602 is to
8 gather information from the physical machines in the P2P system that are represented by
9 the DHT logical space 604. Another exemplary use of the SOMO tree structure 602 is to
10 disseminate information to those physical machines. The information collected can be
11 metadata. Metadata describes information regarding the operation of the P2P system,
12 such as information that reflects the behavior of its physical machines. The information
13 that is disseminated to the physical machines can represent instructions that can govern
14 the operation of the physical machines. One can thus interpret the SOMO mechanism as
15 performing a converge cast from the SOMO leaf nodes to the SOMO root node to
16 provide data gathering, and then performing a multicast back down to the SOMO leaf
17 nodes to provide data dissemination.

18 For example, Fig. 8 represents a scenario 802 in which a SOMO tree
19 structure 804 is being used to collect information from physical machines 806 in
20 the P2P system via the DHT logical space 808. More specifically, the leaf SOMO
21 nodes retrieve the required information from their hosting DHT nodes. (As a side-
22 effect, this procedure can also restart a child SOMO node if it has disappeared
23 because its hosting DHT node crash has crashed). One or more applications 810
24 can invoke this gathering operation for any defined purpose (such as for the
25 purpose of performance monitoring, that is, collecting information regarding

1 various loads and capacities of the physical infrastructure that comprises the P2P
2 system).

3 More specifically, Fig. 8 depicts the configuration of the SOMO tree structure 804
4 to gather information by showing arrows that point from each SOMO node to its
5 corresponding parent SOMO node. In this manner, information funnels up the SOMO
6 tree structure 804 from its leaf SOMO nodes to its root SOMO node. The application(s)
7 810 can extract a complete report from the root SOMO node that culls information from
8 the entire P2P system. This report can contain raw unorganized data. Alternatively, this
9 report can contain merged and sorted data provided that the SOMO nodes have been
10 configured to perform this function before passing the information that they collect onto
11 to their corresponding parent SOMO nodes. The SOMO nodes can be configured to
12 perform this task by configuring the *op* member to perform merging and sorting. This is
13 merely one illustrative example. The SOMO nodes can execute other operations on the
14 information as it passes through the SOMO nodes on its way to the root SOMO node,
15 such as various arithmetic operations.

16 The following pseudo-code provides one technique for gathering information
17 using the SOMO tree structure 804.

18
19 Pseudo-Code Excerpt No. 7: SOMO gathering procedure

```
20 get_report (SOMO_node s) {  
21     Report_type rep[1..k]  
22     for i ∈ [1..k]  
23         if (s.child[i] ≠ NULL) // retrieving via DHT  
24             rep[i] = deref(s.child[i]).report  
25     s.report = s.op(rep[])
```

1 }

2

3 To gather system metadata, the SOMO nodes can periodically perform the above
4 procedure by requesting reports from their respective children.

5 The gather procedure can be tuned to extract specific information from the SOMO
6 tree structure 804. More specifically, the hierarchical nature of the SOMO tree structure
7 804 facilitates the use of complex range queries to discover information relevant to a
8 given logical DHT space region. For example, if k is 2, and it is desired to retrieve a
9 status report of the first quarter of the DHT logical space, an application need only obtain
10 a report from the left child SOMO node of the second-level SOMO tree structure 810
11 (e.g., SOMO child node 812). Another useful implementation involves registering
12 queries at SOMO nodes, which essentially transforms the SOMO mechanism into a
13 publish-subscribe (“pub-sub”) infrastructure.

14 In contrast, Fig. 9 shows a scenario 902 in which a SOMO tree structure 904 is
15 being used to disseminate information to physical machines 906 in the P2P system via the
16 DHT logical space 908. One or more applications 910 can invoke this dissemination
17 operation for any defined purpose (such as for disseminating instructions to the physical
18 machines 906). The configuration of the SOMO tree structure 904 to disseminate
19 information is represented in Fig. 9 by showing arrows that point from parent SOMO
20 nodes to their respective child SOMO nodes. In this manner, information propagates
21 down the SOMO tree structure 904 from its root SOMO node to its leaf SOMO nodes.
22 The information can be propagated through the SOMO tree structure 904 without
23 modification by the SOMO nodes. Alternatively, by virtue of their *op* member, the
24 SOMO nodes can perform any kind of operation on the information before it is passed to
25 their associated child SOMO nodes. Also, as described for the case of data gathering, it

1 is possible to disseminate information to only parts of the DHT logical space 908 by
2 involving only selected parts of the SOMO tree structure 904.

3 In general, the procedures shown in Figs. 8 and 9 describe scalable and fault
4 tolerant mechanisms that can be used in a number of information gathering and
5 disseminating applications. For example, the SOMO mechanism can be used to monitor
6 the health of a P2P system and to provide a UI interface to clients (which alerts these
7 clients to this health information). This can be implemented by gathering information
8 from various performance counters installed on machines used in the P2P system.

9 Additional applications and variations of the data overlay and SOMO tree
10 structure can be implemented. For example, in one exemplary implementation, the
11 SOMO mechanism can be configured to locate multiple classes of critical nodes in the
12 routing structure of the SOMO tree structure; these critical nodes are referred to as
13 supernodes. To provide this application, the report type provided by the SOMO nodes is
14 set to "sorted list," and the *op* provided by the SOMO nodes is set to "merge-sort." The
15 supernodes provide reports that are available at various internal SOMO nodes, where the
16 root SOMO node provides a complete report. These supernodes can accordingly act as
17 indexing or routing hubs.

18 In another exemplary implementation, the SOMO technique can be used to
19 discover a density map of node capacities. Such information can guide object placement
20 in the DHT logical space, or can be used to migrate nodes from dense regions to less
21 dense regions in the DHT logical space.

22 In another exemplary implementation, the SOMO mechanism can be configured
23 to implement a publish-subscribe ("pub-sub") infrastructure.

24 In another exemplary implementation, the SOMO mechanism can create an image
25 of a single resource pool comprised of DHT nodes forming the DHT logic space.

1 In another exemplary implementation, various mechanisms can be used to reduce
2 or eliminate potential harm caused by malicious attack (where users attempt to access
3 system resources for unauthorized and potentially destructive purposes). That is, the
4 SOMO mechanism has been described above principally in a collaborative environment
5 in which malicious attack by users is not expected. But safeguards can be taken if the
6 SOMO tree structure is applied to systems in which malicious DHT nodes (corresponding
7 to malicious user entities) may be present. For example, assume that denial-of-service
8 attacks are mounted by relentlessly requesting the SOMO root node. One way to defend
9 against this attack is by propagating copies throughout the P2P network, thereby stopping
10 the attacks on the edge of the system. Alternatively, assume that there is a potential that
11 the SOMO reports can be compromised in multiple ways on the paths used for
12 information aggregation and/or dissemination. To guard against this, reports can be
13 signed. Also, multiple SOMO internal nodes can generate redundant SOMO reports;
14 various voting schemes can then be used to ascertain consensus and intruder detection.
15 Alternatively, assume that there is a potential that an individual SOMO node can simply
16 “cheat” or “lie” about its own status (e.g., pass inaccurate information about its status).
17 Feedback processes among peers external to the SOMO mechanism can be used to
18 establish trustworthiness of each SOMO node.

19 In another exemplary implementation, increased bandwidth can be provided at
20 SOMO nodes near the top of the SOMO tree structure to accommodate heavy traffic
21 loads in this part of the SOMO tree structure. This bandwidth concern can also be
22 addressed, in part, by only gathering information that reflects a change in the system
23 (e.g., by only gathering delta Δ information). Bandwidth can be further reduced by
24 compressing the information that is transferred through the SOMO tree structure.

1 In another exemplary implementation, functionality is provided for locating a
2 particular node using the SOMO mechanism. An application can then swap this
3 particular SOMO node with the SOMO node which is currently hosting the SOMO root
4 node. This allows the SOMO mechanism to optimize itself.

5 Still more advanced exemplary implementations can be built by providing
6 algorithms that act on the metadata that is gathered from the SOMO tree structure, or
7 which generate the information that is propagated down through the SOMO tree
8 structure. In other words, more advanced techniques can be implemented by providing
9 suitable functionality in the application layers (810, 910) shown in Figs. 8 and 9,
10 respectively.

11 12 *C.3. Performance*

13 The data overlay, and particularly the SOMO tree structure, provides an elegant
14 and flexible mechanism for collecting or disseminating information. The mechanism is
15 flexible in the sense that it specifies neither the type of information it should gather
16 and/or disseminate, nor the operation invoked to process this information. In other
17 words, SOMO operations are programmable. Further, using the abstraction of a data
18 overlay (especially the host routing shortcut), the SOMO tree structure's performance is
19 also insensitive to the particular implementation of the hosting DHT system.

20 In general, both the gathering and dissemination phases are $O(\log_k N)$ bounded,
21 where N is the total number of entities in the P2P system (where the "Big O" notation is
22 used here to express the efficiency of the SOMO algorithm). Each operation in SOMO
23 involves no more than $k+1$ interactions, making it fully distributed. Further, data in the
24 SOMO tree structure can be regenerated in $O(\log_k N)$ time. The SOMO tree self-
25 organizes and self-heals in the same time bound.

More specifically, in one example, assume that the data gathering procedure described above is periodically executed at an interval of T . In this case, information is gathered from the SOMO leaf nodes and this information flows to its SOMO root node with a maximum delay of $\log_k N \cdot T$. This bound reflects the case in which flows between hierarchies are completely unsynchronized. If the upper SOMO nodes' requests immediately trigger the same actions in their children, then the latency can be reduced to $T + t_{hop} \cdot \log_k N$, where t_{hop} is the average latency of a trip in the hosting DHT. The unsynchronized flow has latency bound of $\log_k N \cdot T$, whereas the synchronized version will be bounded by T in practice (e.g., 5 minutes in one exemplary implementation). Note that $O(t_{hop} \cdot \log_k N)$ is the absolute lower bound in one exemplary implementation. For 2M SOMO nodes and with $k=8$ and a typical latency of 200ms per DHT hop, the SOMO root node will have a global view with a lag of 1.6s. The dissemination scenario offers similar performance to that described above for the information gathering scenario.

D. Exemplary Computer Environment for Implementing One P2P Participant

The data overlay described above is a data structure that can be spread out over multiple machines and possibly over other infrastructure in a P2P system. Thus, each of the participants in the P2P system can be viewed as implementing a part of the data overlay. To achieve this effect, each participant can store the necessary code and data to create the data overlay and to interact with it. This code and data can be stored in the volatile and/or non-volatile memory of each participant (to be described below).

For example, Fig. 10 shows a high level view of one exemplary P2P participant 1000. This participant 1000 corresponds to a general purpose computer or server type computer 1002 and an associated display device 1004. However, the participant 1000

1 can be implemented using other kinds of computing equipment. For example, although
2 not shown, the participant 1000 can include hand-held or laptop devices, set top boxes,
3 mainframe computers, etc.

4 Exemplary computer 1002 includes one or more processors or processing units
5 1006, a system memory 1008, and a bus 1010. The bus 1010 connects various system
6 components together. For instance, the bus 1010 connects the processor 1006 to the
7 system memory 1008. The bus 1010 can be implemented using any kind of bus structure
8 or combination of bus structures, including a memory bus or memory controller, a
9 peripheral bus, an accelerated graphics port, and a processor or local bus using any of a
10 variety of bus architectures. For example, such architectures can include an Industry
11 Standard Architecture (ISA) bus, a Micro Channel Architecture (MCA) bus, an Enhanced
12 ISA (EISA) bus, a Video Electronics Standards Association (VESA) local bus, and a
13 Peripheral Component Interconnects (PCI) bus also known as a Mezzanine bus.

14 Computer 1002 can also include a variety of computer readable media, including
15 a variety of types of volatile and non-volatile media, each of which can be removable or
16 non-removable. For example, system memory 1008 includes computer readable media in
17 the form of volatile memory, such as random access memory (RAM) 1012, and non-
18 volatile memory, such as read only memory (ROM) 1014. ROM 1014 includes an
19 input/output system (BIOS) 1016 that contains the basic routines that help to transfer
20 information between elements within computer 1002, such as during start-up. RAM
21 1012 typically contains data and/or program modules in a form that can be quickly
22 accessed by processing unit 1006.

23 Other kinds of computer storage media include a hard disk drive 1018 for reading
24 from and writing to a non-removable, non-volatile magnetic media, a magnetic disk drive
25 1020 for reading from and writing to a removable, non-volatile magnetic disk 1022 (e.g.,

1 a “floppy disk”), and an optical disk drive 1024 for reading from and/or writing to a
2 removable, non-volatile optical disk 1026 such as a CD-ROM, DVD-ROM, or other
3 optical media. The hard disk drive 1018, magnetic disk drive 1020, and optical disk drive
4 1024 are each connected to the system bus 1010 by one or more data media interfaces
5 1028. Alternatively, the hard disk drive 1018, magnetic disk drive 1020, and optical disk
6 drive 1024 can be connected to the system bus 1010 by a SCSI interface (not shown), or
7 other coupling mechanism. Although not shown, the computer 1002 can include other
8 types of computer readable media, such as magnetic cassettes or other magnetic storage
9 devices, flash memory cards, CD-ROM, digital versatile disks (DVD) or other optical
10 storage, electrically erasable programmable read-only memory (EEPROM), etc.

11 Generally, the above-identified computer readable media provide non-volatile
12 storage of computer readable instructions, data structures, program modules, and other
13 data for use by computer 1002. For instance, the readable media can store the operating
14 system 1030, application programs 1032, other program modules 1034, and program data
15 1036.

16 The participant 1000 can include a variety of input devices. For instance, the
17 participant 1000 includes the keyboard 1038 and a pointing device 1040 (e.g., a “mouse”)
18 for entering commands and information into computer 1002. The participant 1000 can
19 include other input devices (not illustrated), such as a microphone, joystick, game pad,
20 satellite dish, serial port, scanner, card reading devices, digital or video camera, etc.
21 Input/output interfaces 1042 couple the input devices to the processing unit 1006. More
22 generally, input devices can be coupled to the computer 1002 through any kind of
23 interface and bus structures, such as a parallel port, serial port, game port, universal serial
24 bus (USB) port, etc.
25

1 The participant 1000 also includes the display device 1004. A video adapter 1044
2 couples the display device 1004 to the bus 1010. In addition to the display device 1004,
3 the participant 1000 can include other output peripheral devices, such as speakers (not
4 shown), a printer (not shown), etc.

5 Computer 1002 operates in a peer-to-peer networked environment using logical
6 connections to one or more remote computers, such as a remote computing device 1046.
7 The remote computing device 1046 can comprise any kind of computer equipment,
8 including another general purpose personal computer, portable computer, a server, etc.
9 Remote computing device 1046 can include all of the features discussed above with
10 respect to computer 1002, or some subset thereof.

11 Any type of network 1048 can be used to couple the computer 1002 with remote
12 computing device 1046, such as the WAN 402 of Fig. 4, an intranet, a LAN, etc. The
13 computer 1002 couples to the network 1048 via network interface 1050, which can utilize
14 broadband connectivity, modem connectivity, DSL connectivity, or other connection
15 strategy. Although not illustrated, the participant 1000 can provide wireless
16 communication functionality for connecting computer 1002 with remote computing
17 device 1046 (e.g., via modulated radio signals, modulated infrared signals, etc.).

18 Any of the above-identified stores in the computer 1002 can be used to store the
19 code and data used to implement part of a data overlay, such as part of the SOMO tree
20 structure.

21
22 To conclude, a strategy was described for building a data structure on top of a
23 DHT in a P2P system. A specific hierarchical tree structure was specifically described
24 for disseminating information into the P2P system and for collecting information from
25 the P2P system.

1 Certain operations were described as constituting distinct steps performed in a
2 certain order. Such implementations are exemplary and non-limiting. Certain steps
3 described herein can be grouped together and performed in a single operation, and certain
4 steps can be performed in an order that differs from the order employed in the examples
5 set forth in this disclosure.

6 Further, a number of examples will be presented in this disclosure in the
7 alternative (e.g., case A or case B). In addition, this disclosure encompasses those cases
8 which combine alternatives in a single implementation (e.g., case A and case B), even
9 though this disclosure may not expressly mention these conjunctive cases in every
10 instance.

11 Generally, although the invention has been described in language specific to
12 structural features and/or methodological acts, it is to be understood that the invention
13 defined in the appended claims is not necessarily limited to the specific features or acts
14 described. Rather, the specific features and acts are disclosed as exemplary forms of
15 implementing the claimed invention.